

# Diffusion Models

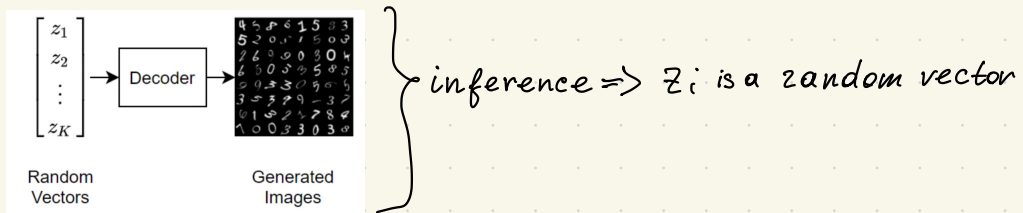
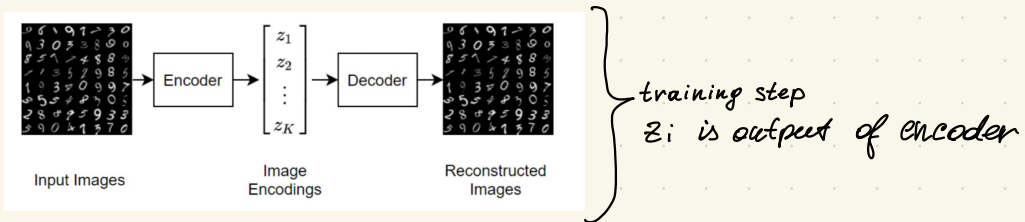
- ① Variational Autoencoders
- ② Diffusion Models as particular case of VAE
- ③ SDE
- ④ LD and FP eq.
- ⑤ DM as score-matching
- ⑥ Classifier guidance

## ① Variational autoencoder (VAE)

VAE is a neural network that learns to produce its input

Given objects  $X = (x_1, \dots, x_n)$ , where  $x_i \in \mathbb{R}^D$ , construct  $z_i \in \mathbb{R}^d$ , with which we can reconstruct  $X$ .

### 1.1. Architecture overview



12 Training objective: maximize the likelihood of  $X$ ;  $\log p(X|\theta) \rightarrow \max$   
For the VAE model:  $p(X, z|\theta)$

$$\log p(X|\theta) = \int p(X, z|\theta) dz \quad \left. \begin{array}{l} \text{this integral is intractable} \\ \text{because relationship between } X \text{ and } z \\ \text{is highly non-linear} \end{array} \right\}$$

$\int p(x, z | \theta) dz$  is intractable, that's why **variational L-bound is used**:

$$\log p(x | \theta) \geq \mathcal{L}(q, \theta) = \int q(z | x, \psi) \log \frac{p(x, z | \theta)}{q(z | x, \psi)} dz \rightarrow \max_{\psi, \theta}$$

$\{\psi, \theta\}$  = parameters

**variational inference; ELBO**

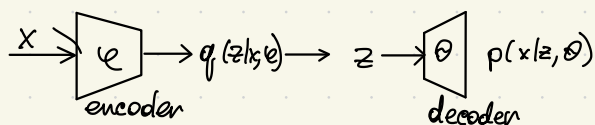
$$\log p(x | \theta) \geq \underbrace{\int q(z | x, \psi) \log p(x | z, \theta) dz}_{\text{Reconstruction error}} -$$

$$\underbrace{\int q(z | x, \psi) \log \frac{q(z | x, \psi)}{p(z)} dz}_{\text{KL-divergence (Regularizer)}}$$

Reconstruction error

KL-divergence  
(Regularizer)

1.3 VAE scheme with  $\psi, \theta$

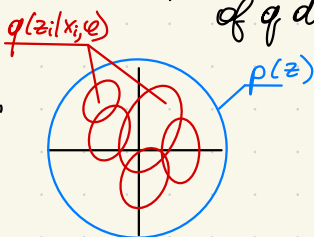


1.4 meaning of loss components

- Why use **regularizer**?

To avoid **point** distributions; with **KL** term we penalize large deviations of  $q$  distributions from  $z$

With KL term,  
we learn  $q$ s  
close to  $p(z)$ !



- Without the **reconstruction error**, we would have similar  $q$ s.

1.5 Disadvantages of VAE

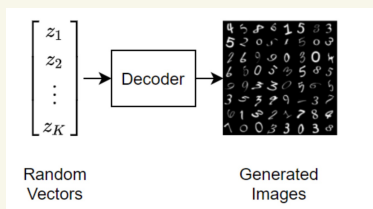
$p(z) \neq \frac{1}{n} \sum q_j(z_j | x_j, \psi) \rightarrow$  We cannot cover all  $p(z)$  with  $q$ s.

$p(z) \neq \frac{1}{n} \sum_j q_j(z_j | x_j, \phi) \rightarrow$  in generative mode, we take  $z$  from prior distribution:

$$\hat{z} \sim p(z)$$

$$x \sim p(x | \hat{z}, \phi)$$

(sampling)



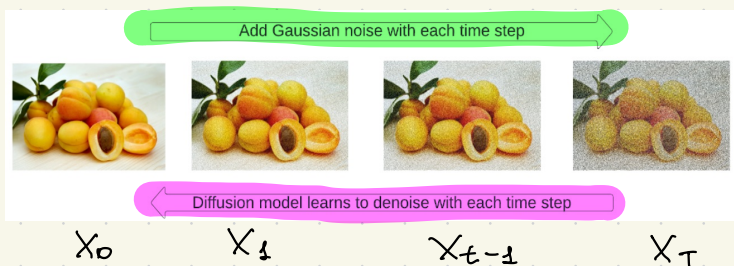
That's why in generative mode we do not have realistic output.

VAEs are good for reconstruction, but not good at generating high-quality images

Exception: if  $d > D$  (Very deep VAE), we can generate better images

## ② Diffusion Models as particular case of VAEs

DMs are also latent variable generative models (like VAEs), but they work by gradually adding noise to the input data:



2.1 Model overview

Diffusion model is defined by:

- forward
- backward processes
- sampling

2.2

## Diffusion Model modes

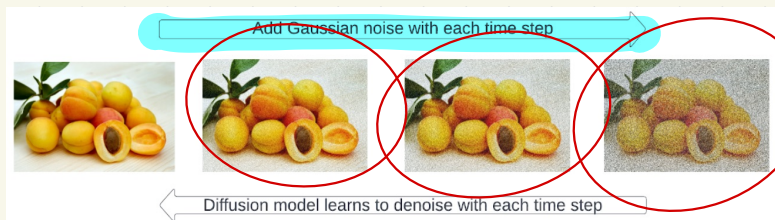
Training

Add noise (forward)  
Reconstruct from noised image (backward)

Inference

Sample from  $z$   
Reconstruct from noisy image

## 2.3 Forward process



$x_0, \dots, x_t \rightarrow$  representations of input  
At each  $t$  step, we add noise

$$x_{t+1} = \sqrt{1-\beta} x_t + \sqrt{\beta} \epsilon \quad \beta \ll 1 \quad \epsilon \sim \mathcal{N}(\epsilon | 0, I)$$

$\downarrow$  constant  $\downarrow$  noise

in general:  $x_t = \sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \epsilon$ , where  $\alpha_t = (1-\beta)^t$

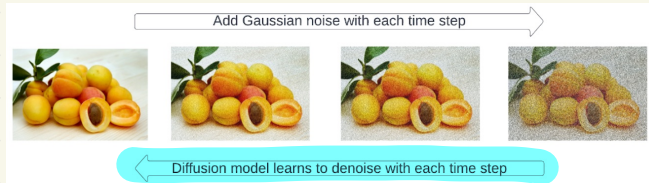
$$\alpha_T \ll 1, \alpha_T \approx 0 \quad q(x_T | x_0) \approx \mathcal{N}(x_T | 0, I)$$

$$\rightarrow q(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\alpha_t} x_0, (1-\alpha_t)) \xrightarrow{t \rightarrow T} \mathcal{N}(x_T | 0, I)$$

$$q(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\alpha_t} x_0, (1 - \alpha_t) \mathbf{I}) \xrightarrow{t \rightarrow T} \mathcal{N}(x_T | 0, \mathbf{I})$$

Any object from training distribution diffuses to  $\mathcal{N}$ !  
So, DMs do not have VAE's drawback

## 2.4 Backward process (Training objective)



$$\log p(\text{Obs} | \theta) \geq \int q(\text{Mid} | \text{Obs}, \epsilon) \log \frac{p(\text{Obs}, \text{Mid} | \theta)}{q(\text{Mid} | \text{Obs}, \epsilon)}$$

VAEs { Observed = Obs  
Hidden = Mid

Training objective  
(see 1.2)

$$\log p(x_0 | \theta) \geq \int \underbrace{q(x_1 \dots x_T | x_0)}_{\text{we don't have } \epsilon, \text{ } q \text{ is fixed}} \log \frac{p(x_0, x_1 \dots x_T | \theta)}{q(x_1 \dots x_T | x_0)} dx_1 \dots dx_T \quad \textcircled{=}$$

we don't have  $\epsilon$ ,  
 $q$  is fixed:  $x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon$

$\textcircled{=}$  { apply product rule }  $\textcircled{=}$

$$\textcircled{=} \int q(x_1 \dots x_T | x_0) \log p(x_0 | x_1, \theta) dx_1 \dots dx_T + \int q(x_1 \dots x_T | x_0) \log \frac{p(x_1 \dots x_T | \theta)}{q(x_1 \dots x_T | x_0)} dx_1 \dots dx_T$$

$$= \int q(x_1 | x_0) \log p(x_0 | x_1, \theta) dx_1 + \int q(x_1 \dots x_T | x_0) \log \frac{p(x_T) p(x_{T-1} | x_T, \theta) \dots p(x_1 | x_2, \theta)}{q(x_T | x_0) q(x_{T-1} | x_T, x_0) \dots q(x_1 | x_2, x_0)} dx_1 \dots dx_T$$

$- \text{KL}(q(x_1 \dots x_T) || p(x_1 \dots x_T))$

$$\int q(x_1|x_0) \log p(x_0|x_1, \theta) dx_1 + \int q(x_1 \dots x_T|x_0) \log \left( \frac{p(x_T) p(x_{T-1}|x_T, \theta) \dots p(x_1|x_2, \theta)}{q(x_T|x_0) q(x_{T-1}|x_T, x_0) \dots q(x_1|x_2, x_0)} \right) dx_1 \dots dx_T =$$

KLS

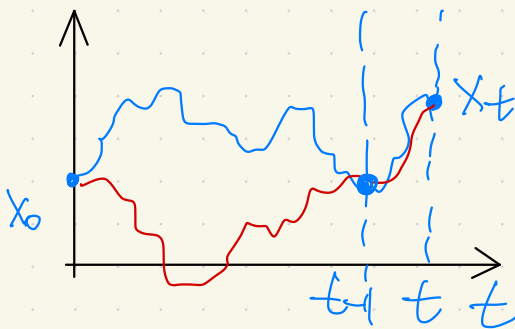
$$= \int q(x_1|x_0) \log p(x_0|x_1, \theta) dx_1 - \sum_{t=2}^T \int q(x_t|x_0) \text{KL}(q(x_{t-1}|x_t, x_0) \| p(x_{t-1}|x_t, \theta)) - \text{KL}(q(x_T|x_0) \| p(x_T)) \rightarrow \max_{\theta}$$

That is why DMs are efficient!

## Training objective

$$\log p(x_0|\theta) \geq \mathcal{L}(\theta) \approx - \sum_{t=2}^T \mathbb{E}_{x_t} \text{KL}(q(x_{t-1}|x_t, x_0) \| p(x_{t-1}|x_t, \theta)) \rightarrow \max_{\theta}$$

$$\begin{aligned} q(x_{t-1}|x_t, x_0) &= \frac{q(x_t|x_{t-1}, x_0) q(x_{t-1}|x_0)}{q(x_t|x_0)} = \\ &= \frac{q(x_t|x_{t-1}) q(x_{t-1}|x_0)}{q(x_t|x_0)} = \\ &= \mathcal{N}(x_{t-1} | \mu(x_0, x_t), \Sigma_{t-1}) \end{aligned}$$



$$\mu(x_0, x_t) = \frac{\sqrt{\alpha_t} \beta}{1 - \alpha_t} x_0 + \frac{\sqrt{1 - \beta} (1 - \alpha_{t-1})}{1 - \alpha_t} x_t$$

$$\tilde{\beta}_t = \frac{1 - \alpha_{t-1}}{1 - \alpha_t} \cdot \beta$$

$$q(x_{t-1} | x_t, x_0)$$

$$p(x_{t-1} | x_t, \theta) = q(x_{t-1} | x_t, \underbrace{x_\theta(x_t, t)}}_{\text{trained with deep NN}})$$

trained with deep NN

Then, we can rewrite training objective:

$$\text{KL}(q(x_{t-1} | x_t, x_0) \| p(x_{t-1} | x_t, x_0)) = \text{const} \cdot \|x_0 - x_\theta(x_t, t)\|_2$$

## 2.5 Training procedure

- ① Take  $x_0$  from dataset
- ② Take arbitrary  $\tau \in [2, T]$
- ③ Generate  $x_\tau \sim q(x_\tau | x_0)$   $x_\tau = \sqrt{\alpha_\tau} x_0 + \sqrt{1 - \alpha_\tau} \epsilon$
- ④ Differentiate  $\text{KL}(q(x_{\tau-1} | x_0) \| p(x_{\tau-1} | x_\tau, \theta))$  w.r.t.  $\theta$   
 $\text{const} \cdot \|x_0 - x_\theta(x_\tau, \tau)\|^2$